



# Measuring and Quantifying Interestingness of Documents

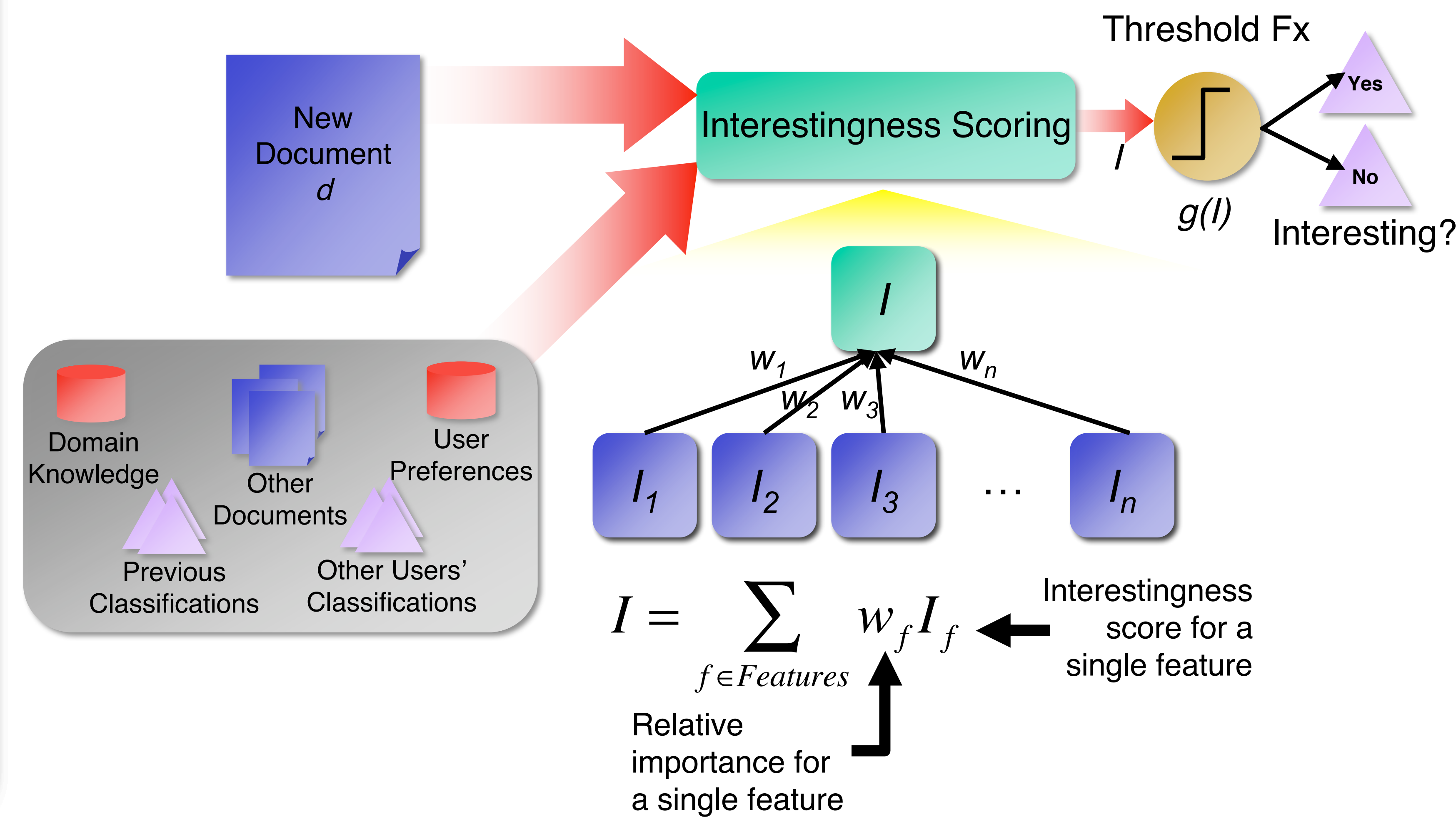


Raymond K. Pon, UC Los Angeles

David Buttler and Terence Critchlow, Center for Applied Scientific Computing  
Lawrence Livermore National Laboratory

Search engines, such as Google, assign scores to documents, such as news articles, based on their relevancy to a particular query or topic. However, not all relevant documents may be interesting to the user, particularly if the document is old or yields very little new information. Furthermore, relevancy scores do not take into account the uniqueness and the impact to a domain of interest by the document. To remedy this situation, we present a general framework for defining and measuring the “interestingness” of documents taking into account the above document features along with the event the document describes, document source analysis, user-feedback, and collaborative scoring. Applications include intelligent news aggregators, automated event detection and monitoring, and personalized web portals.

## Document Classification

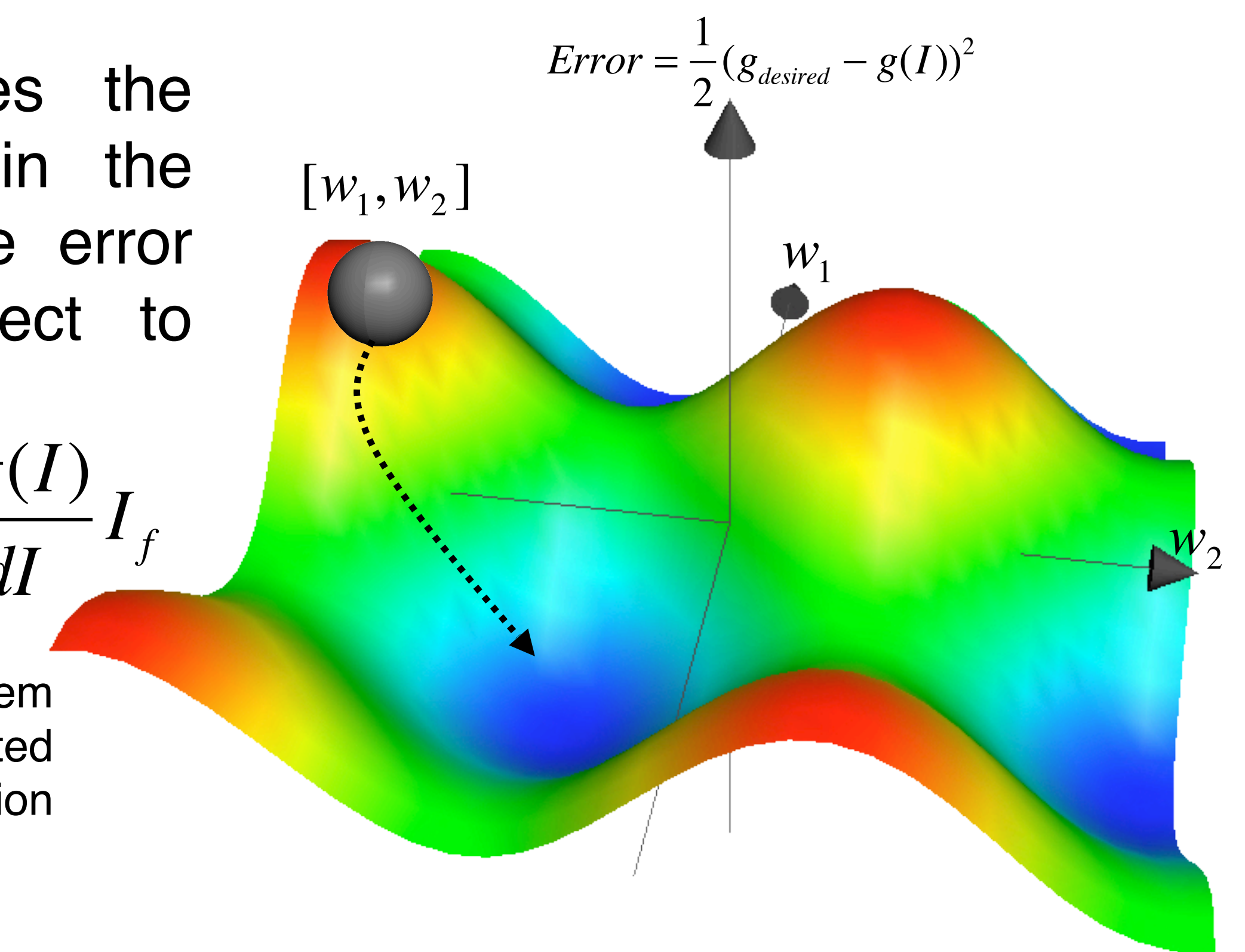


## Training with Gradient Descent

Gradient descent moves the weight assignment  $w_f$  in the opposite direction of the error change rate with respect to weight assignments

$$\Delta w_f = \eta (g_{\text{desired}} - g(I)) \frac{dg(I)}{dI} I_f$$

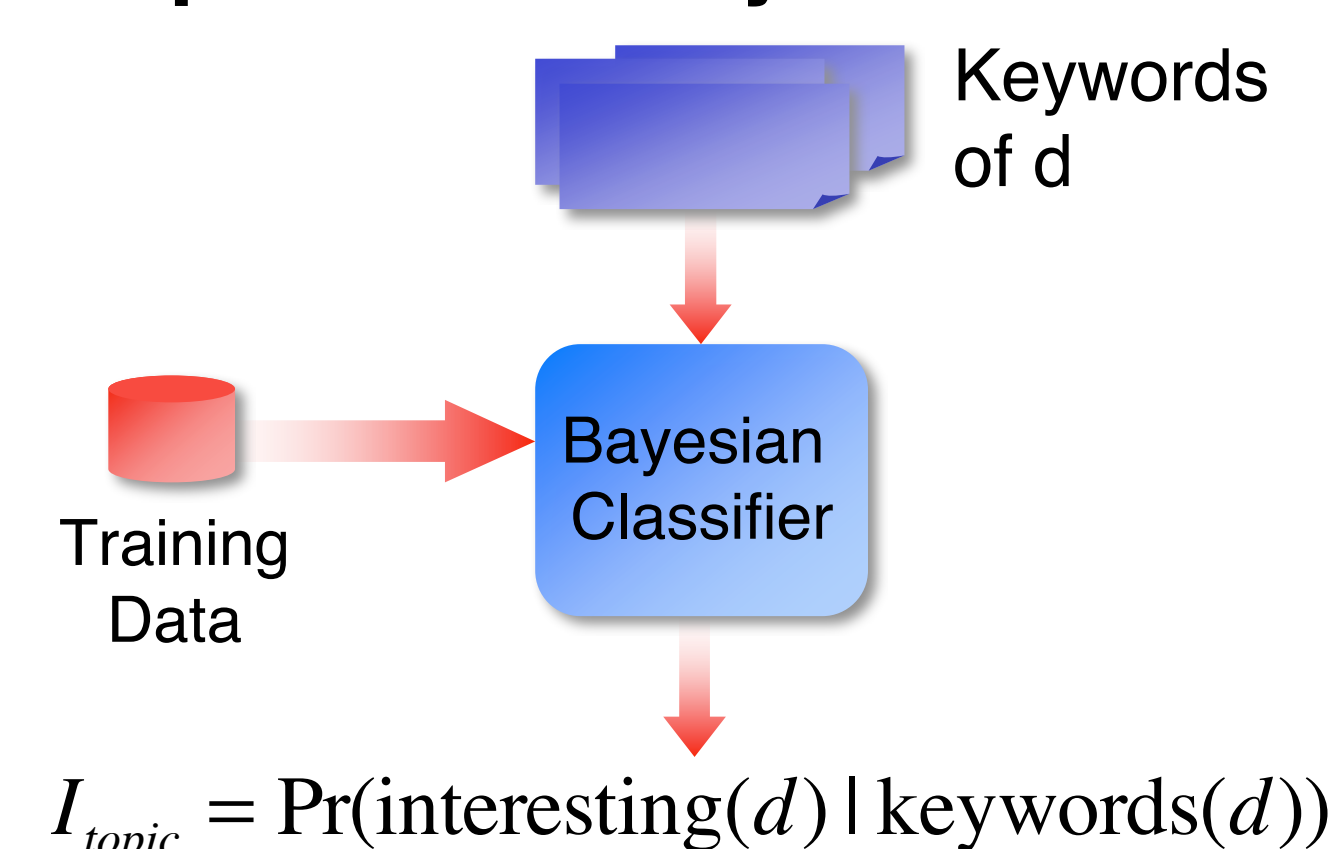
Learning rate (step value)  $\uparrow$  Correct classification  
System generated classification  $\uparrow$



Moving weight assignments for two features to a region with lower error.

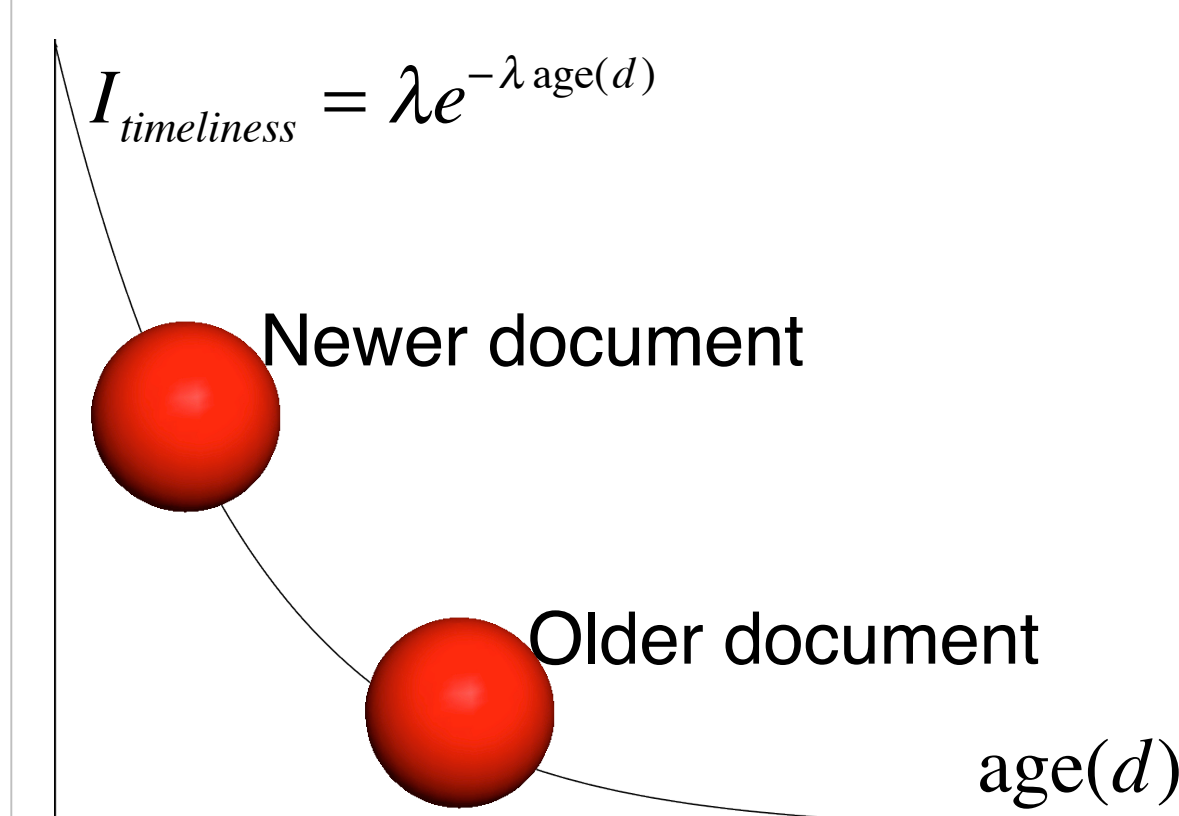
## What Makes a Document Interesting?

### Topic Relevancy



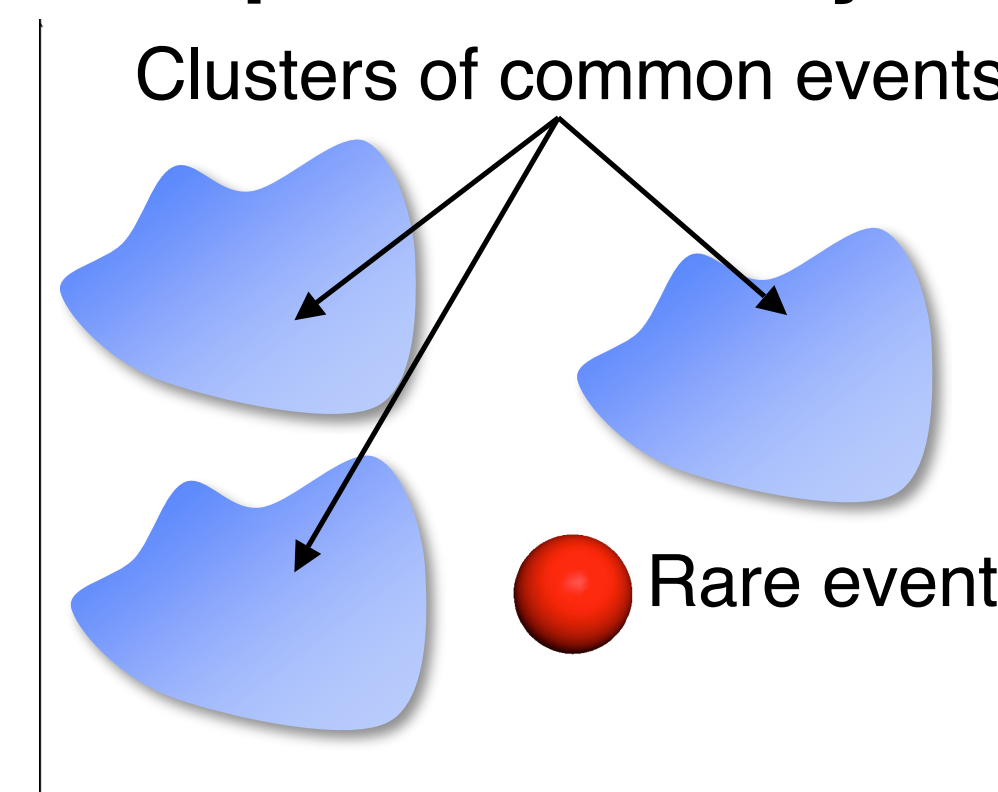
A document must be relevant to a topic of interest for it to be interesting.

### Timeliness



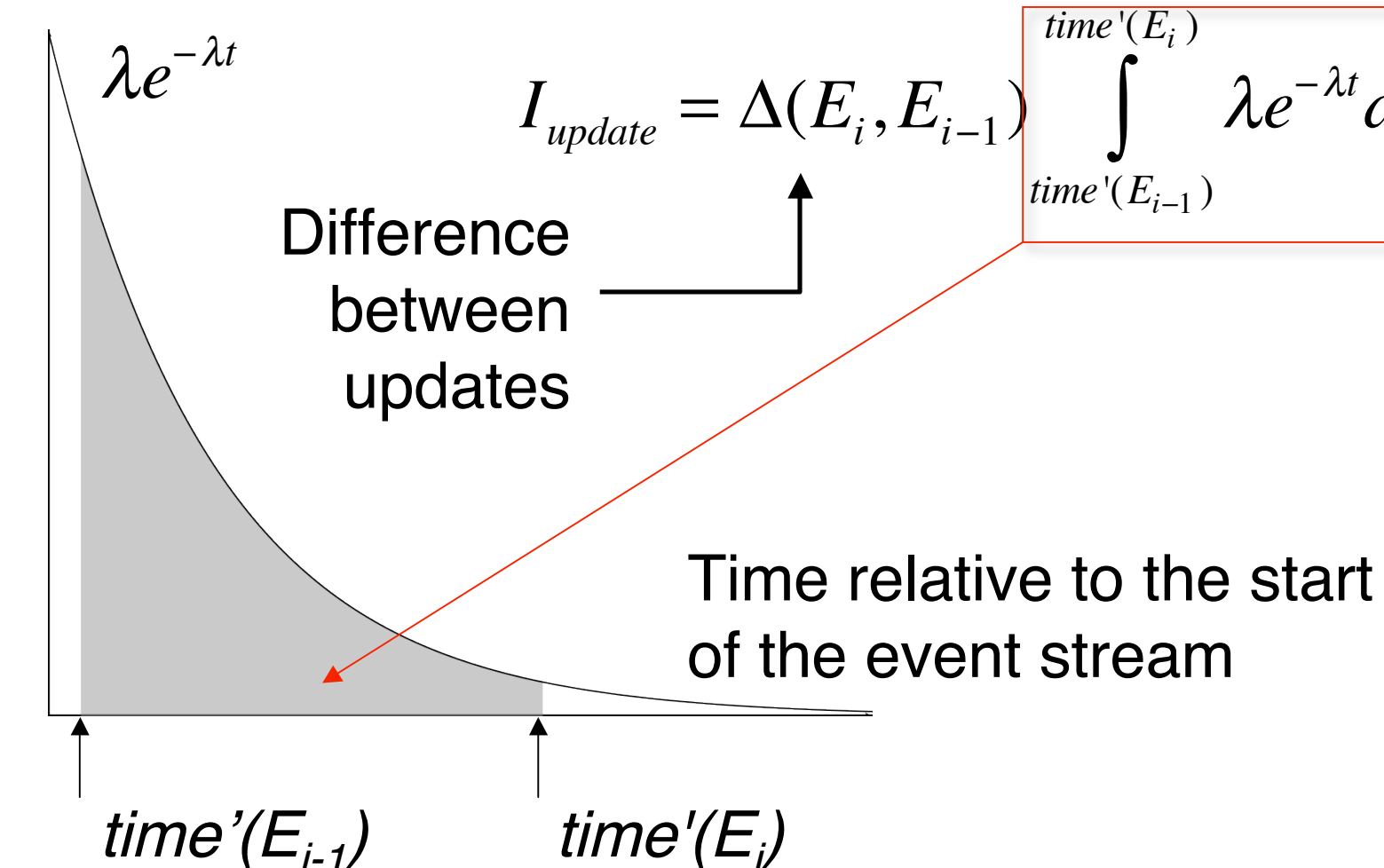
Newer documents are generally more interesting than older documents.

### Uniqueness/Rarity of Event



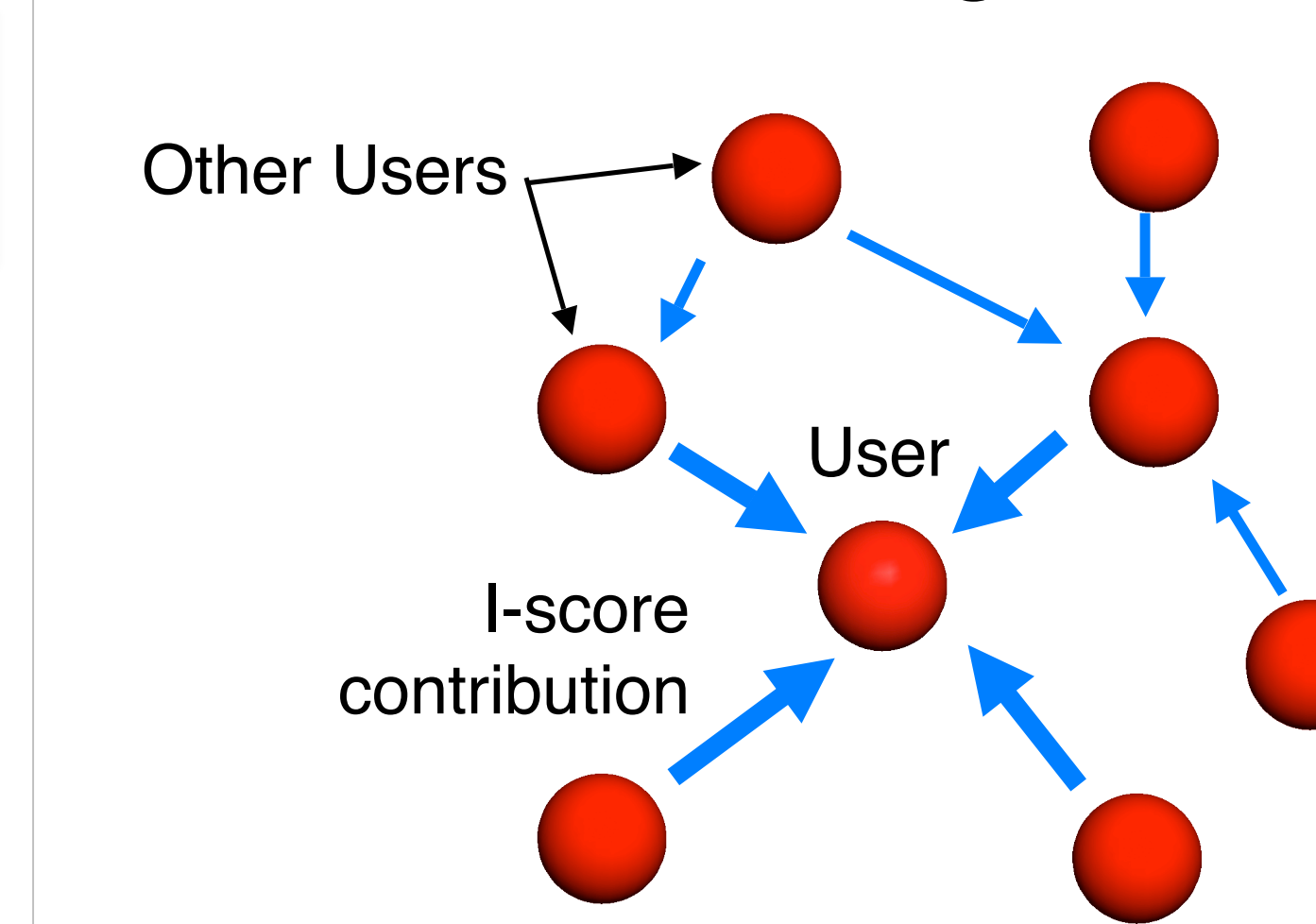
A document describing a rare or unique event is interesting.

### Significance of Event Updates



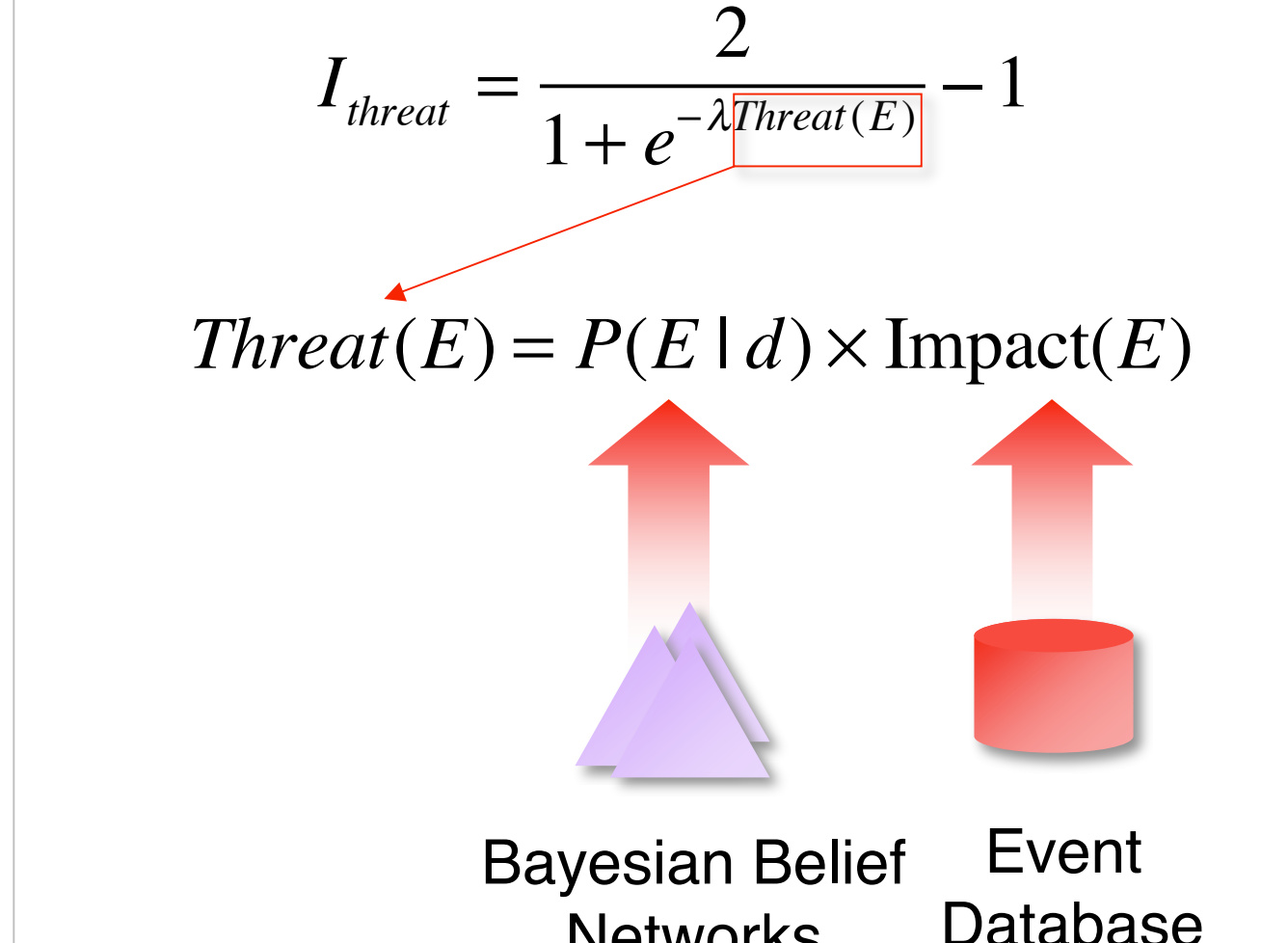
A document providing significant new information to an event last updated a long time ago is interesting. Additionally, documents that first describe a new event are interesting.

### Collaborative Scoring



A document classified as interesting to one user should also be interesting to other users with similar interests. Additionally, documents that are similar to each other should have similar interestingness scores for the same user. Social network analysis can help address this feature.

### Threat Assessment



In the case of documents describing disaster events, such as disease outbreaks, those that describe highly probable events with significant impact are interesting.